# BNL-Lake

## A step towards US-Lake

# Data lake in the US:
# a lot of communication over the last year

# Lake: a 2-D concept

- **Horizontal** dimension: Spatial distribution of storages (the different locations)

- **Vertical** dimension: Hierarchy of storage at a given location

  - Can be very simple (cache only) to complex with a variety of solutions (cache, disk, tape) implementing different technologies

- The 2-D are managed by :

  - An internal file catalog

  - An active information system

- The 2 dimensions can be independently implemented —> **BNL-Lake**

# Caching and lake: A note in passing

- Ilija performed extensives studies and simulations  (Slides)  of caches on US sites

- Ultimate conclusions cannot be reached because of tight connection between data placement (DDM) and workflow management (WFMS)

- Current data organization cannot be used for asserting the benefit of a caching system in a data lake

- Only a real prototype could be used to evaluate the full benefit of a lake

## Last episode summary

Production input are slightly more cacheable (52% accesses and 67% data volume) than Analysis inputs (35% accesses and 37% data volume).

Different file types have very different access patterns (eg. HITS, EVNT, payload files are very cacheable, DAODs, panda*, AODs less so).

Claim: even a cache of 50TB per site would be sufficient to deliver roughly half or the accesses and data volume.

3

## Conclusions II

Smallish caches at all sites could deliver ~ half of data volume.
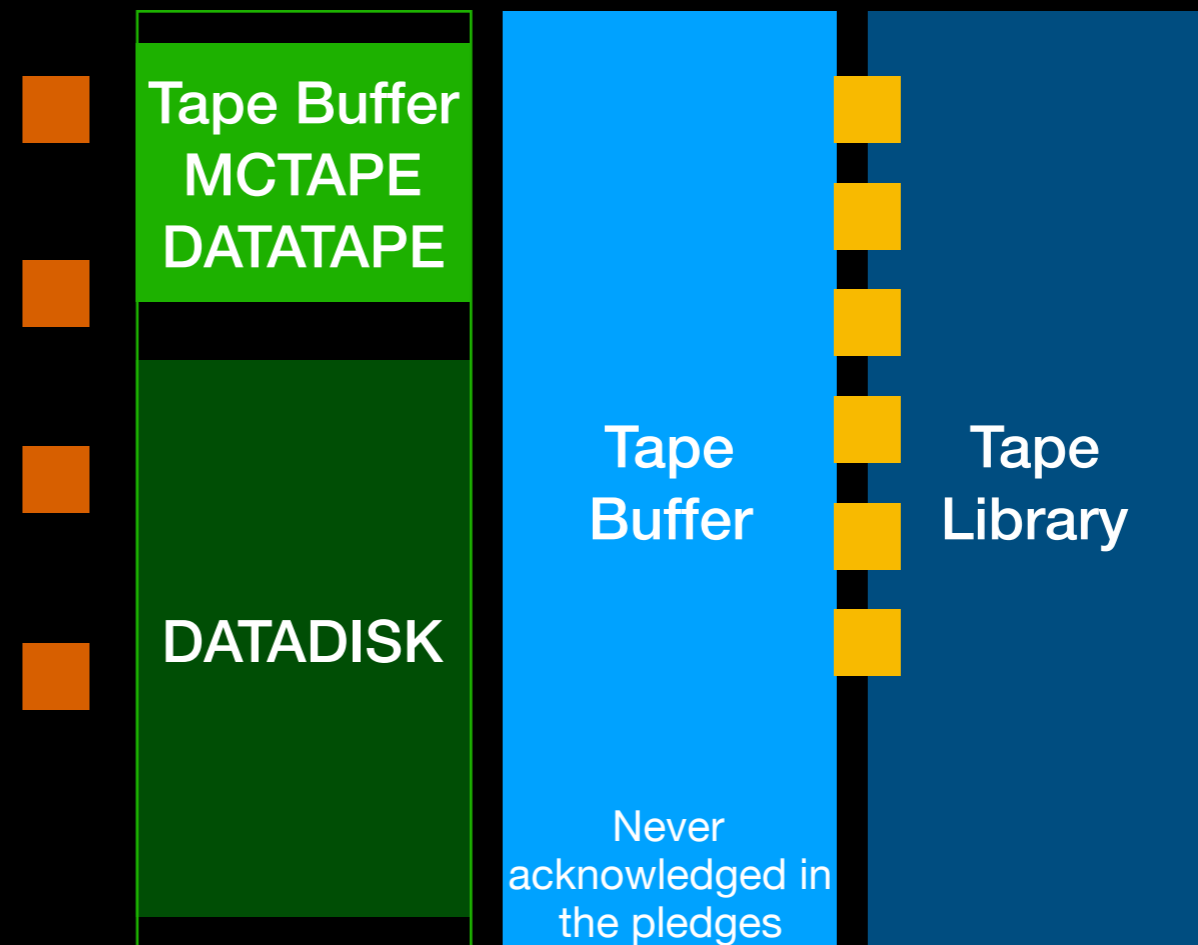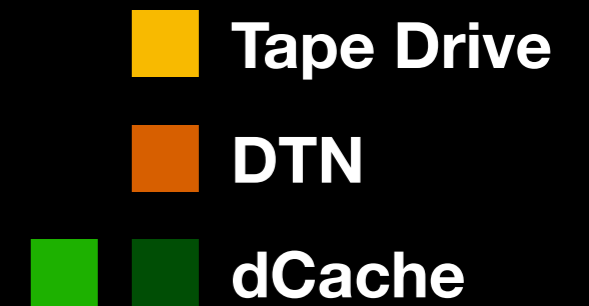
Not much benefit from 2nd level cache.

NEW QUESTION

   If we increase cache size, where would be best to add it?

ANSWER

   At sites
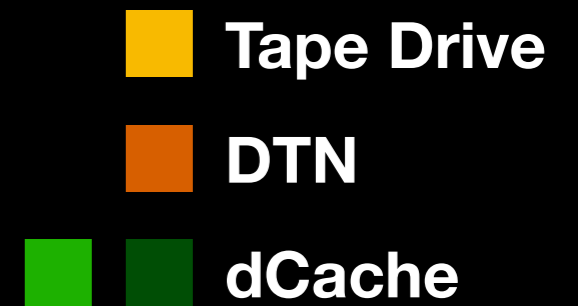
15

# Current setup

**Tape Drive** (yellow)
**DTN** (orange)
**dCache** (green)

Tape Buffer
MCTAPE
DATATAPE

DATADISK

Tape
Buffer

Never
acknowledged in
the pledges

Tape
Library

**None of the internals are specified
by existing MoUs**
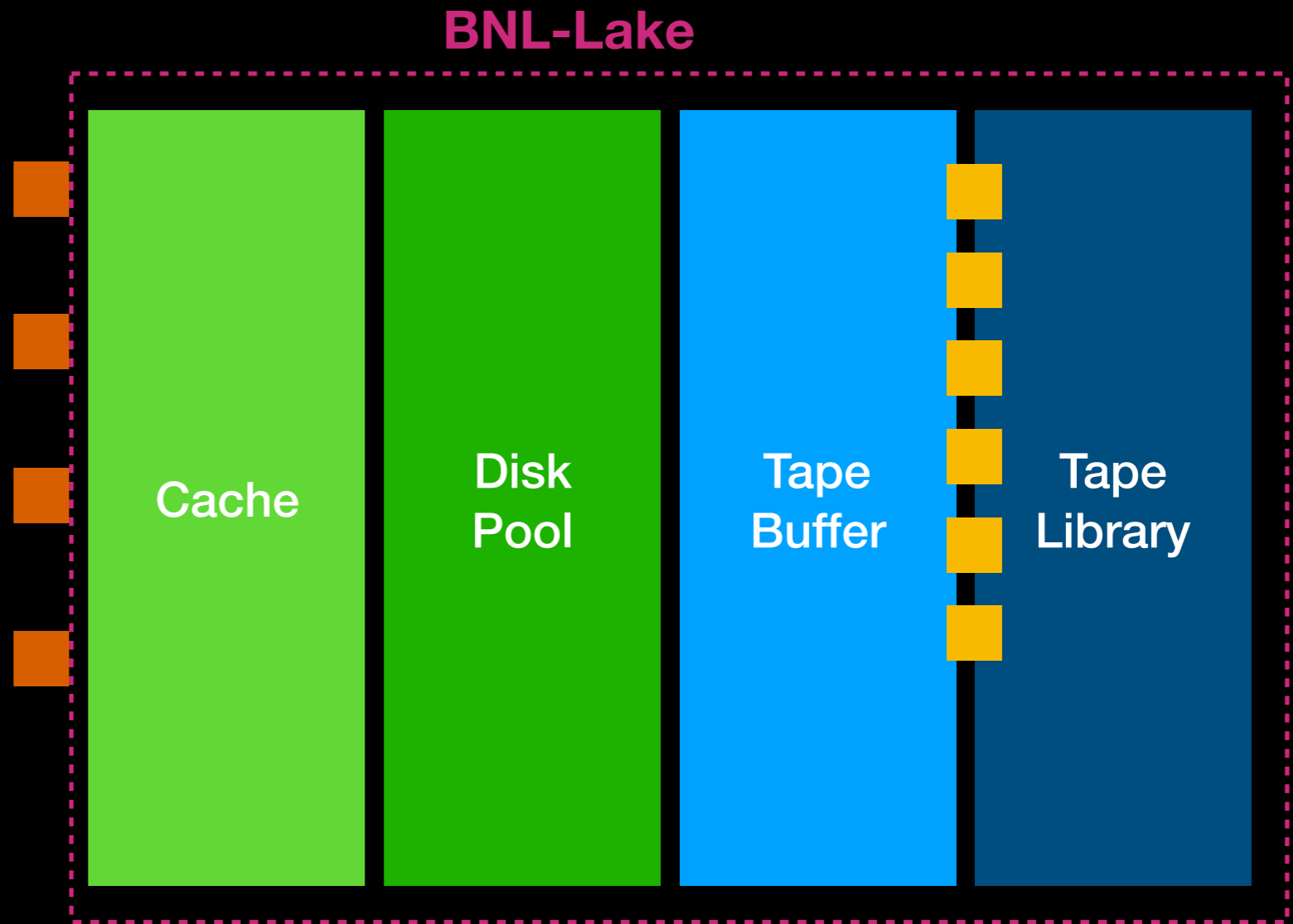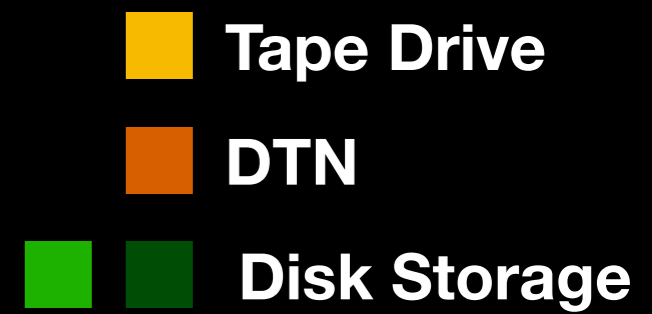
# Current setup

Tape Drive
DTN
dCache

- WFMS/DDM interact with disk tokens and SRM

- The details of site internals are unknown to them (and will remain)

  - Number and type of tape drives

  - Tape buffer size

  - LAN Bandwidth

  - I/O capability and reliability of disk storage,

  - etc… etc… etc…

- Optimisation of data access and storage for given requirements can only be performed by the site

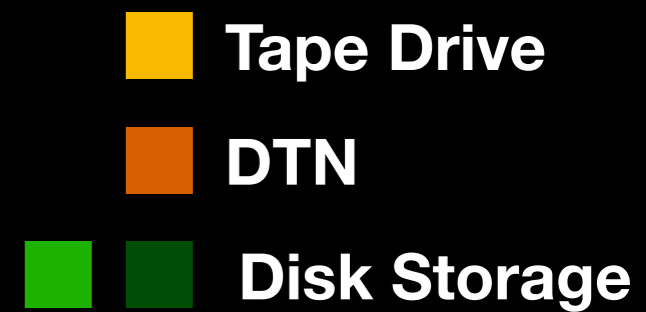Tape Buffer
MCTAPE
DATATAPE

DATADISK

Tape Buffer

Never acknowledged in the pledges

Tape Library

**None of the internals are specified by existing MoUs**

6

# Proposal

# Proposal

**Tape Drive**

**DTN**

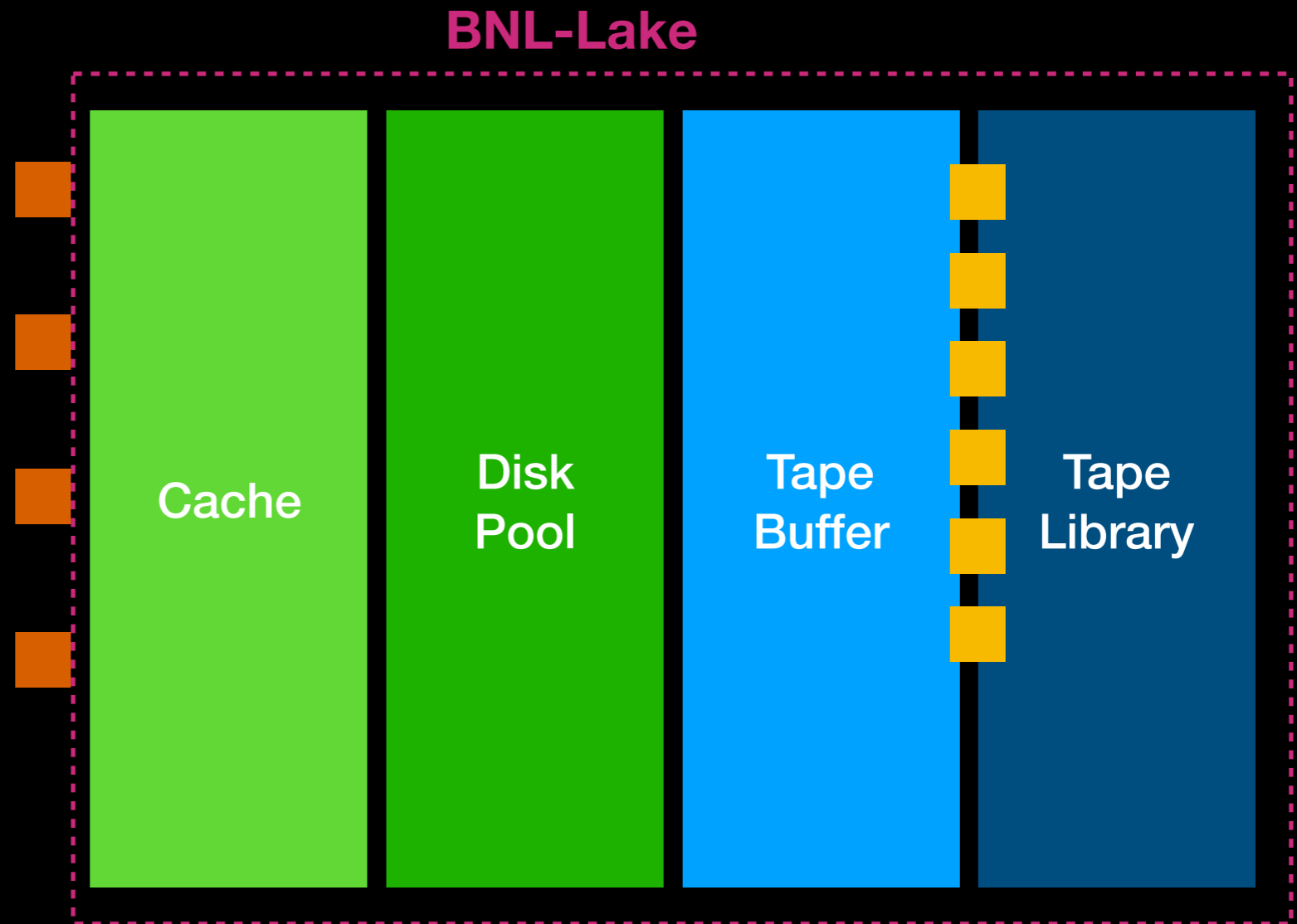**Disk Storage**

**BNL-Lake**

- Can be implemented in several steps in parallel to present infrastructure

  1. One single disk token

  2. Cache

  3. Information system

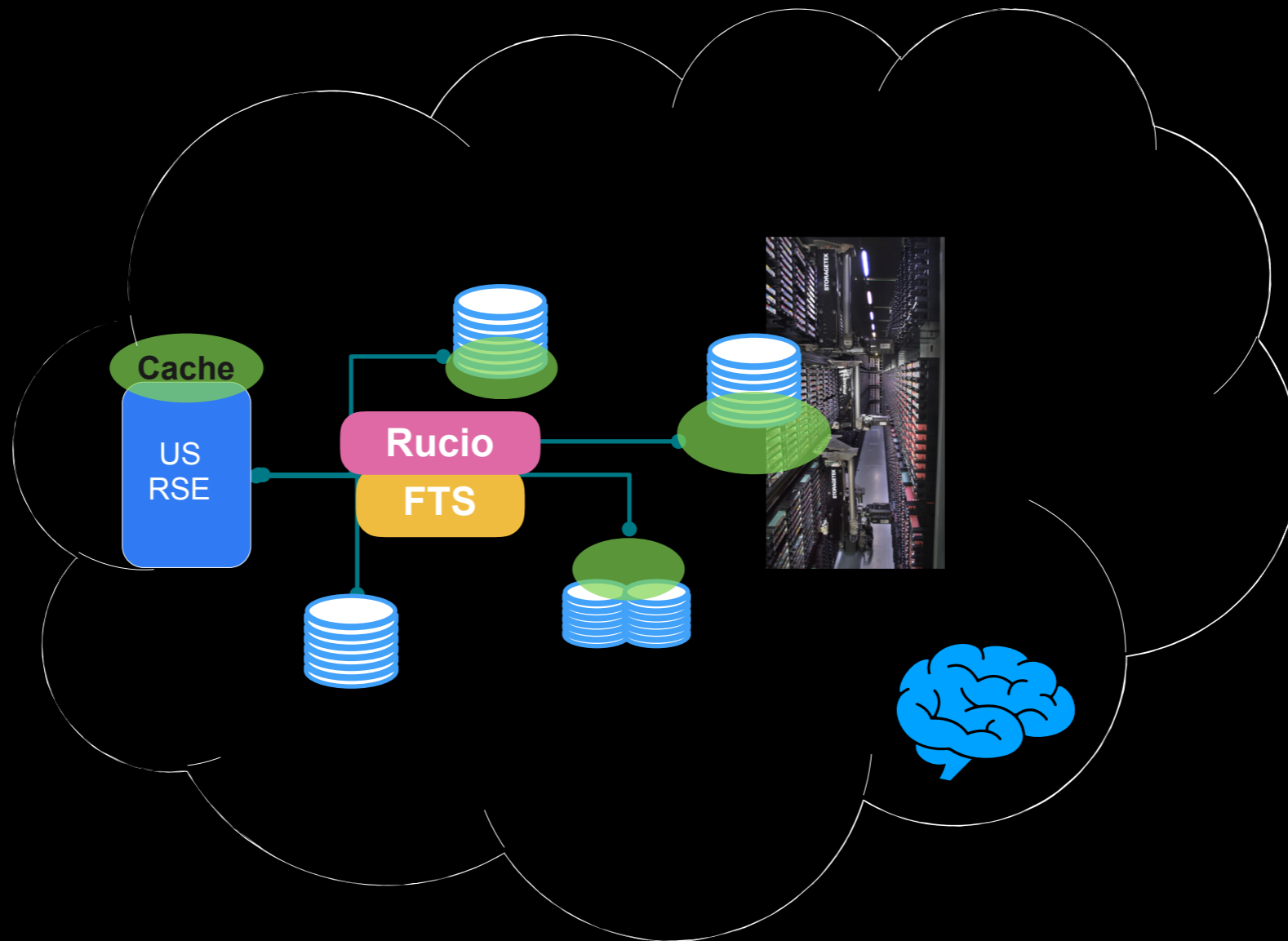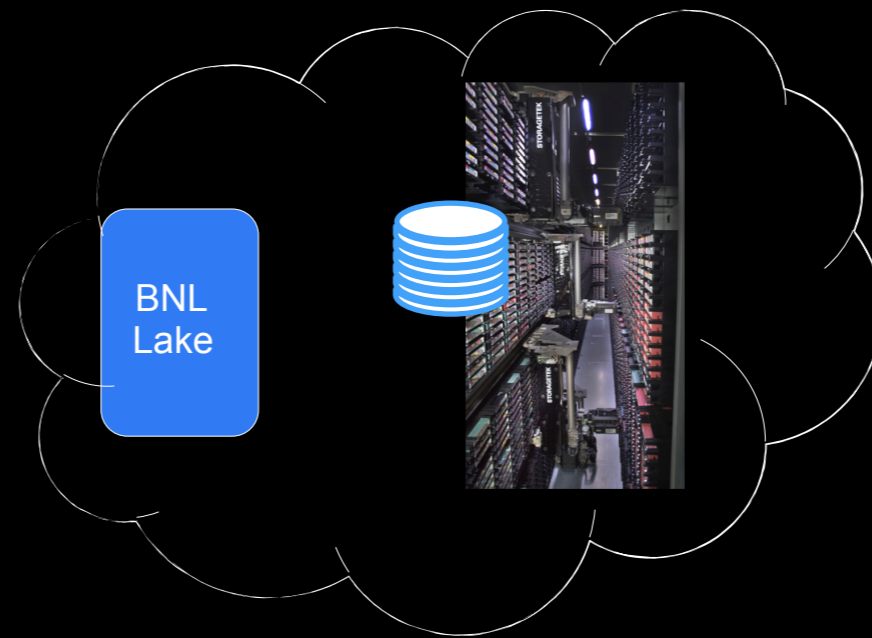- The simplest way to asses the capabilities of the Lake concept without interference of the WFMS/DDM

Cache

Disk Pool

Tape Buffer

Tape Library

**Note: this concept will break the historical WLCG requirement of a given disk space at a given site**
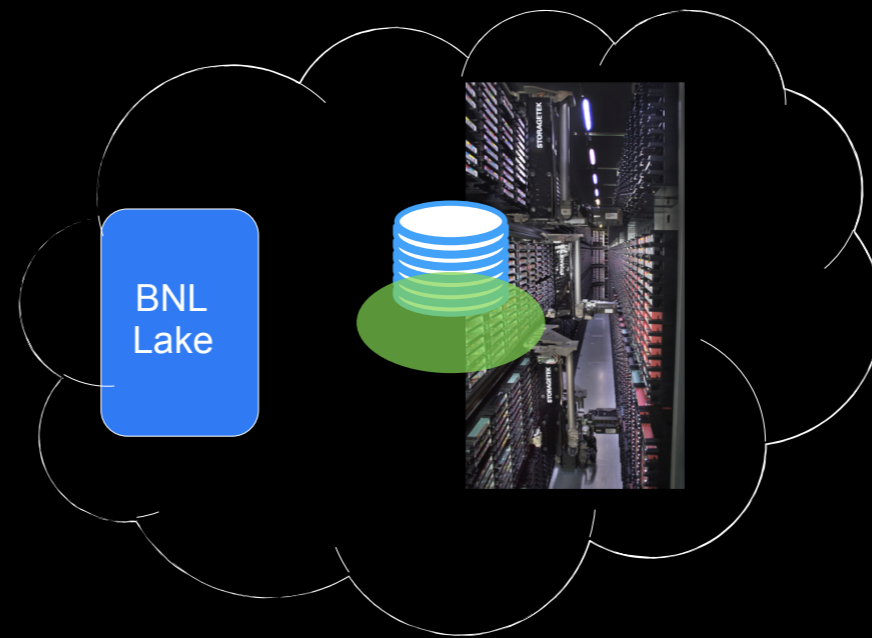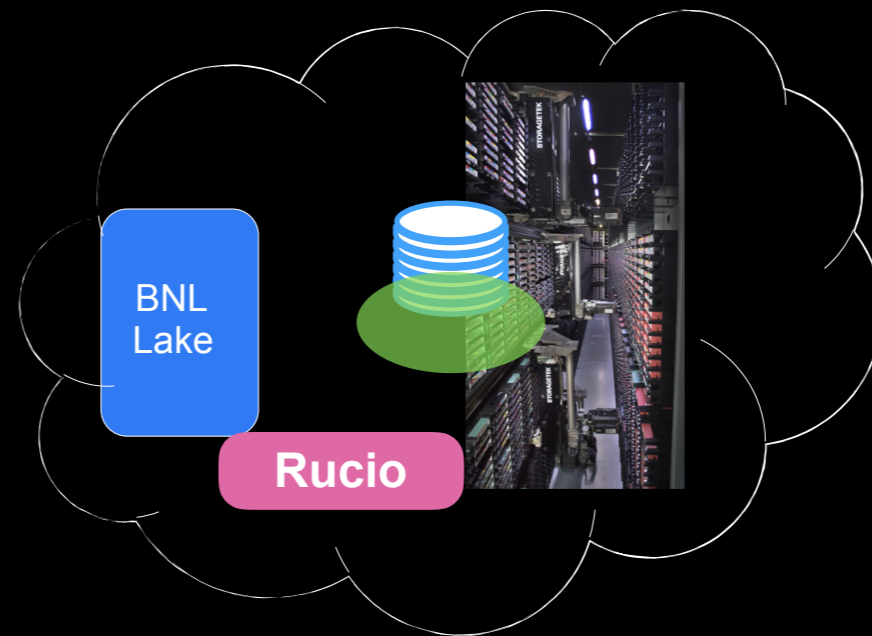
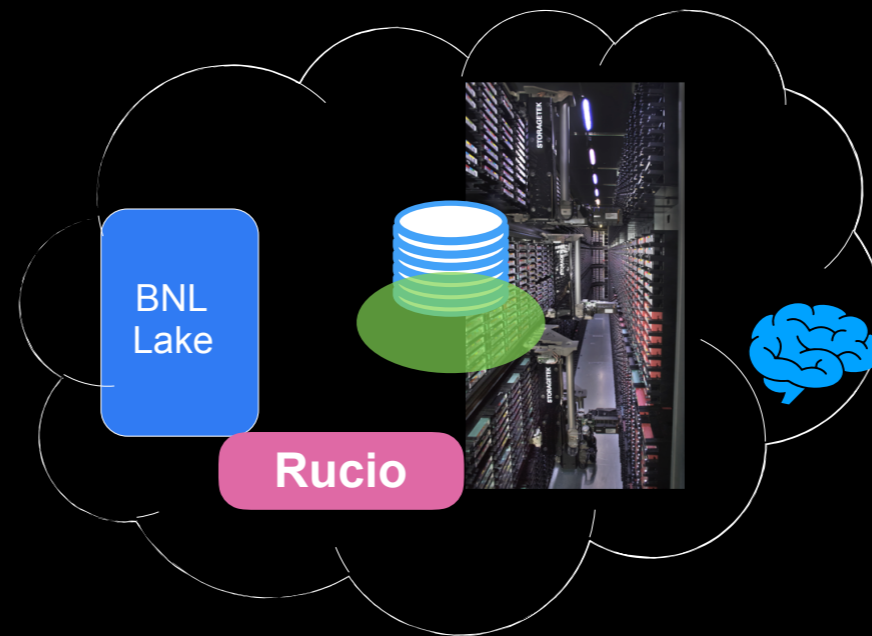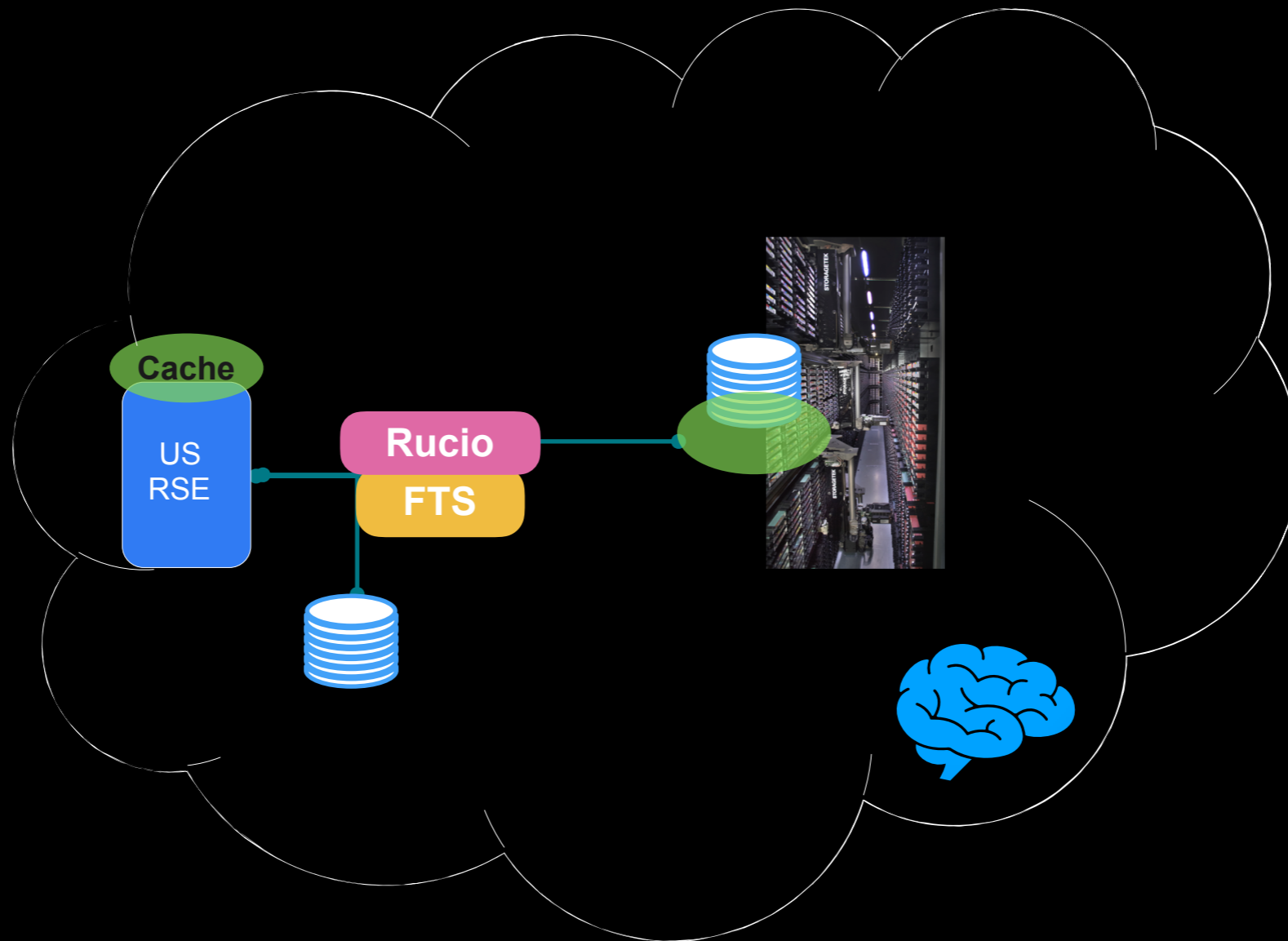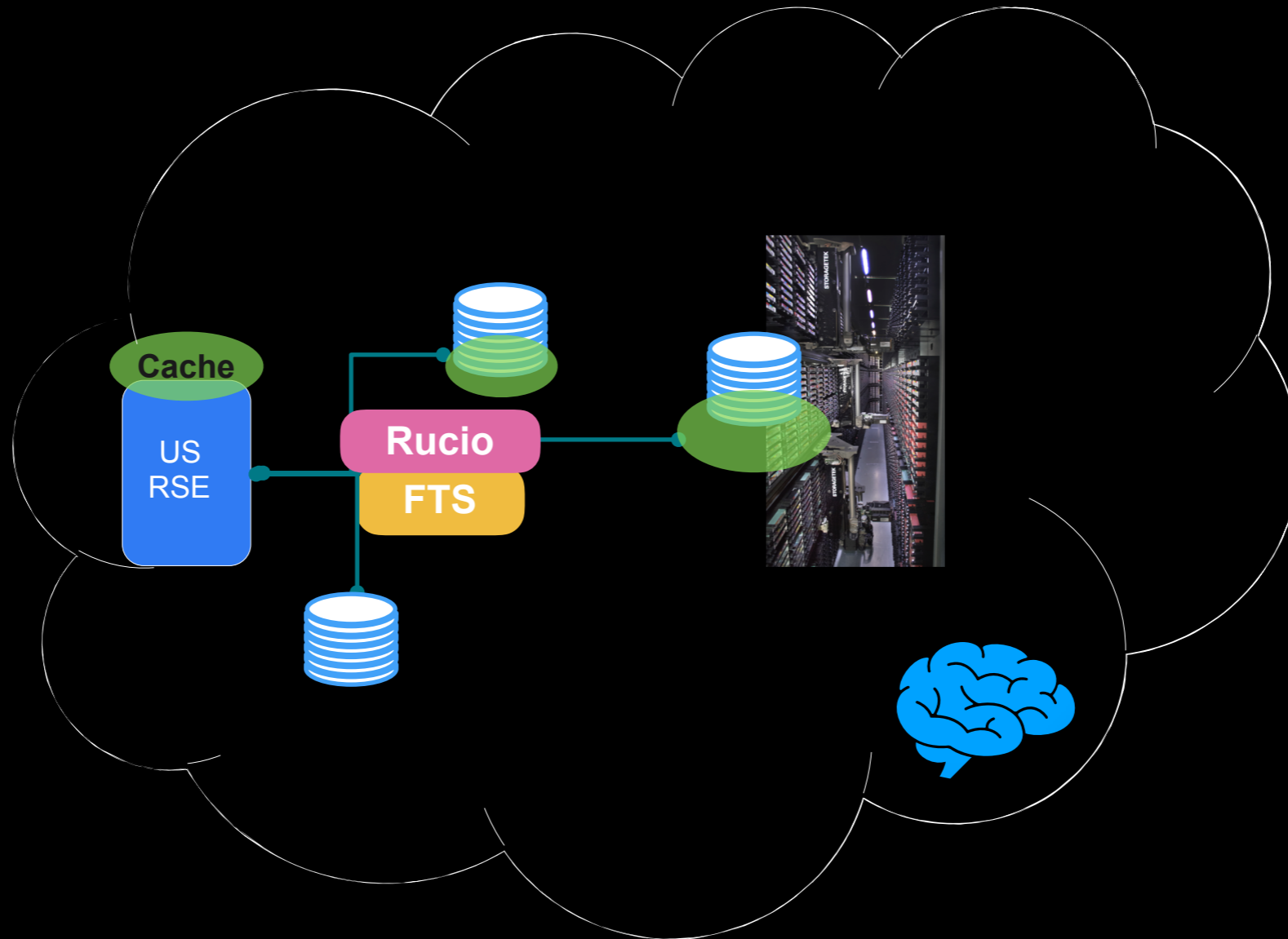# The dream…US-Lake

# PathForward



BNL
Lake

# PathForward

# PathForward

# PathForward

# PathForward

# PathForward

# PathForward